

DATE: 14-2-2017

SUBJECT: معلوماتية

Big data  $\Rightarrow$  refers to exponential growth and availability of data

"Volume" large size - more complex tools - Variety - Volatility

Value or not "القيمة" - Trusted or not "الدقة" - accuracy of data

data structure "Structured - Semi structured - unstructured"

relational  
dB

XML  
JSON

Binary JSON

video  
audio  
pdf

أجزاء  
من  
البيانات

required tools:

• way of storage

• analysis tools

data scientist  $\rightarrow$  Personal Skills "مهارات شخصية"

$\rightarrow$  Technical Knowledge

$\rightarrow$  mathematics tools

OLAP  $\Rightarrow$  online analytical processing

"البيانات التحليلية عبر الإنترنت"

Slide 17 1-not secure

1. easy and simple

Island

2. replication

seperable

3. hacking

"spread sheets"

4-

data ware

house

الخزائن

مخزن البيانات

آخر المطاف

1. Controlled

2. must have permission  
for apply any action

on data

3. Controlled by data Base

administrator

SandBox

permission من قبل

data warehouse

و بعد هذا يتاح للجزء الى جميع البيانات من

على البيانات الخاصة

(((ALAKSA)))



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

وبذلك <sup>for management</sup> High performance

Case study. Bank From local int global

- 1- القوانين الخاصة بكل دولة
- 2- حركه التعامل مع البنك
- 3- High security
- 4- تغير العملاء
- 5- زيادة العملاء فتحتاج الى عالم زيادة
- 6- distributed data

Data scientist - variety

- 2- stream large
- 3- management data
- 4-

Data analyst

يتوقف data محتاجه اليه وليا عما انما توكل لوجده

KPI  $\Rightarrow$  Keep performance indicator.

Classifier  $\rightarrow$  TP, FN, TN, FP  
performance evaluation.

\* Business intelligence وهو يقاس بها

ويجعل analysis على انه اقدر اخذ قرار بناء على هذا التحليل

من يقدر يتعامل الا مع structured data فقط



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

• Data Scientist.

الترجيح والمقارنة السابقة

« unstructured - structured » data مع أي نوع من البيانات

Comparison between them. Scales - Job of each of them

33

Throughput → انت أي شيء يتكلف  
data من شيء لا شيء

Scales of data scientist

lec2

1. Statistics skills

2. Database

+ Adaptive

3. Critical thinking, creative, Innovate + Communications skills

4. machine learning + Data mining + Advanced mathematics.

5. collect data from different online source

6. Extract data & Analysis

7. <sup>can make</sup> Correlations & Connections.

8. web development & web design.

9. programming skills.

85

Data enables. ~ data collectors.

Professional Business reports, machine learning limit



Quantitative

أي يقيس الأشياء بالآلة ويحول  
Technical report  
Statistics view

Skeptical

يتكافح أن يكون مثالي

Communications & Collaborative Skills  
أدائه، ومع التعامل مع الناس Skills

5.10

الأمثلة من مسكن

أولاً pandemics

5.12

Swine Flu

Outcome

بداية الوقاية المرض في مراحل بتوقع

5.13 Life science

Genome

A1. Complexity

A2.

## Life cycle of DAP

lec 3

1. Define problem
2. I/P
3. Issue
4. Outcomes
5. Resources available

2. collect available data and be sure that it's enough and secure.

3. modeling.

4. test model

يفضل بعد الانتهاء من كل مرحلة التأكد أنه تم تسجيلها في ملف  
ملاحظات

documentation.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

جميع الامور حله يتم لتقاربنا لاسر - فانه صبرا الآخر بيم بعد ذلك .  
• additional analysis اي في تحليل  
والقائد validation للبيانات الخاطيه

Business user → end user.

S.9

① Discovery.

- learn about the problem domain.
- hold history of this domain and analysis it
- Documentation of the old project.
- measure to what resource available and period of time and quality of data.

2 - Data Pre  
is responsible of building "Sand Box".

3. model - Sw that will uses.

- 2 - Feature important // Feature selection.
- 3 - HW.

S.16  
probe. استجوب

S.20 discussion in next lecture.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

Big data  $\Rightarrow$  is a popular term which refers to the exponential growth and availability of data, both structured and unstructured.

Three 'V's' to describe the definition of big data.

volume

velocity

variety

Volume.

$\Rightarrow$  There has been a large increase of data volume.

There reasons:

1. All of the transactional data that has added up over the years.
2. Streaming data From social media.
3. machine to machine data increase.

velocity

$\Rightarrow$  Data is being streamed at huge speeds and needs to be dealt with in a timely manner.

1. Social media.

2. mobile devices.

The biggest challenge is how to react fast enough to the massive amount of data that is being flew rapidly.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

## variety

managing all the different formats is an issue many organizations have to battle.

• there are many different types of data

- Structured. • Email. • Audio & video.
- Application data. • Financial transactions.
- Unstructured documents.

• So to manage many organizations have to battle

Volume, velocity, variety, value, veracity.

## Big Data

is data whose scale, distribution, diversity, and/or timeliness require the use of technical architectures and analytics to enable.

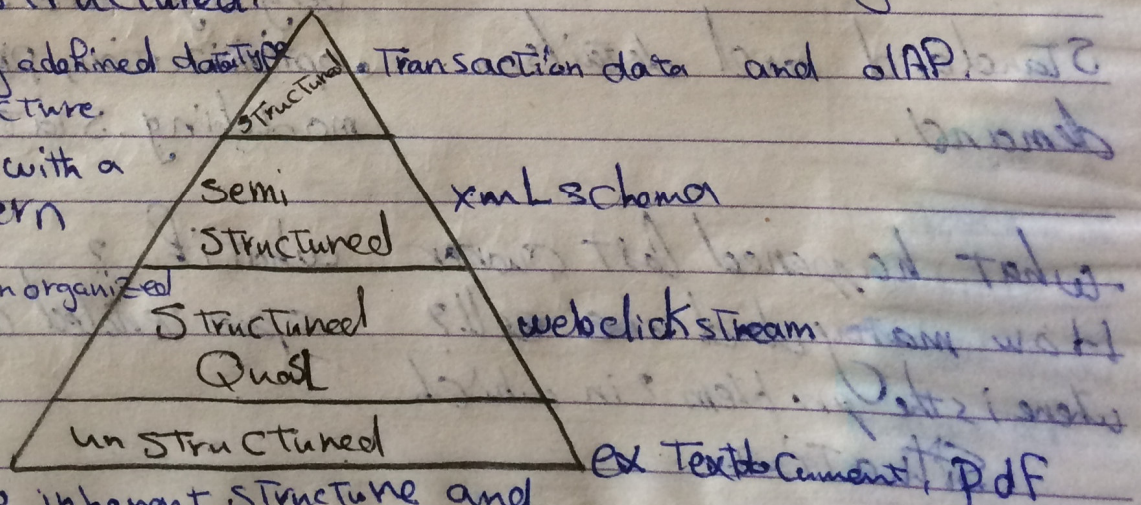
## Key characteristics of BD

- 1 - data volume.
- 2 - processing complexity (parallel computing environment and massively parallel processing)
- 3 - Data Structured.

Data containing defined data type / format, structure.

Textual data files with a discernable pattern

Textual data and unorganized data format



Data has no inherent structure and

is usually stored as different types of files.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

## Storing data

Data islands spreadmarts isolated data	Data warehouse Centralized Data Containers in a purpose built space.	Analytic SandBox Data assets gathered From multiple Sources and the technology For analysis.
• spread sheet and Low volume DB. • Analyst dependent on data extracts.	• Analyst dependent on IT @ DBAs For data access and schema change. • Analysts must spend significant time to get extract For multiple sources.	• enable high performance • reduce cost associated with data replicated. analyst "tunnel" • more robust analyses.

## Business intelligence.

## Data science.

Structured data, traditional  
sources. manageable datasets.

• Structure/unstructured data,  
many types of sources,  
very large data sets.

Standard and details on  
demand.

• optimization, predictive  
modeling. Statistical analysis.

what happened last quarter  
How many did we sell?  
where is the problem & in which  
situation

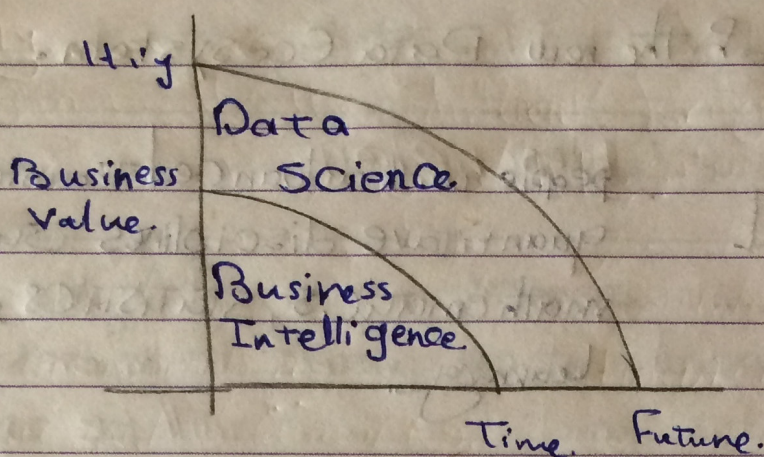
what if...?

open ended questions.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_



implications of typical Architecture For data science.

1. high-value data is hard to reach and leverage
  2. Data is moving in batches From EDW to local analytics.
  3. isolated and analytic projects, rather than centrally-managed or analytics.
- The Big data trend is generating an enormous amount of information that requires advanced analytics and new market players to take advantage of it.

Criteria For Big Data projects.

1. speed of decision making
2. Through put
3. Analysis Flexibility.



DATE: \_\_\_\_\_

SUBJECT: lec 2

## Three Key roles of the new Data Ecosystem.

Data science

- Deep analytical Talent

people with advanced training in quantitative disciplines such as mathematics, statistics, machine learning.

Analysts  
Data savvy managers.

- Data savvy professionals.

people with a basic knowledge of statistics and/or machine learning who can define key questions that can be answered using advanced analytics.

- Technology & Data Enablers.

people providing technical expertise to support analytics projects skill sets including computer programming and DB administrator.

## Data Scientist Key Activities

1. reFrame business challenges as analytics challenges
2. Design implement and deploy statistical models and data mining techniques on big data.
3. Create insights that lead to actionable recommendations.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

## «Data Analytics Life cycle»

• value of using the data Analytics lifecycle.

1. Ensure rigidity and completeness.

2. Enable better transition to members of the cross-Function analytic teams.

Creating and documenting a process will help demonstrate rigor in your findings.

• repeatable.

• Scale to additional analysts.

• Support validity of findings.

Need for a process to Guide Data Science projects.

1. well-defined processes can help guide any analytic project

2. Focus of Data Analytics project Lifecycle is on Data Science projects, not business intelligence.

3. Data science projects tend to require a more consultative approach, and differ from BI projects in a few ways.

• less predictable data

• more projects which lack shape or structure.

• more due diligence in discovery phase.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

## Key roles for a successful Analytic project

Role	Description
Business user	Someone who benefits from the end results and can consult and advise project team on value of end results and how these will be operationalized
project sponsor	person responsible for the genesis of the project, providing the motives for the project and core business problem, generally provide the funding and will measure the degree of value from the final outputs of the working team.
project manager	Ensure key objectives are met on time and at expected quality.
Business Intelligence Analyst	Business domain expertise with deep understanding of the data KPIs, Key metrics and business intelligence from a reporting perspective.
Data Engineer	Deep technical skills to assist with tuning SQL queries for data management extraction and support data realize to analytics sand box



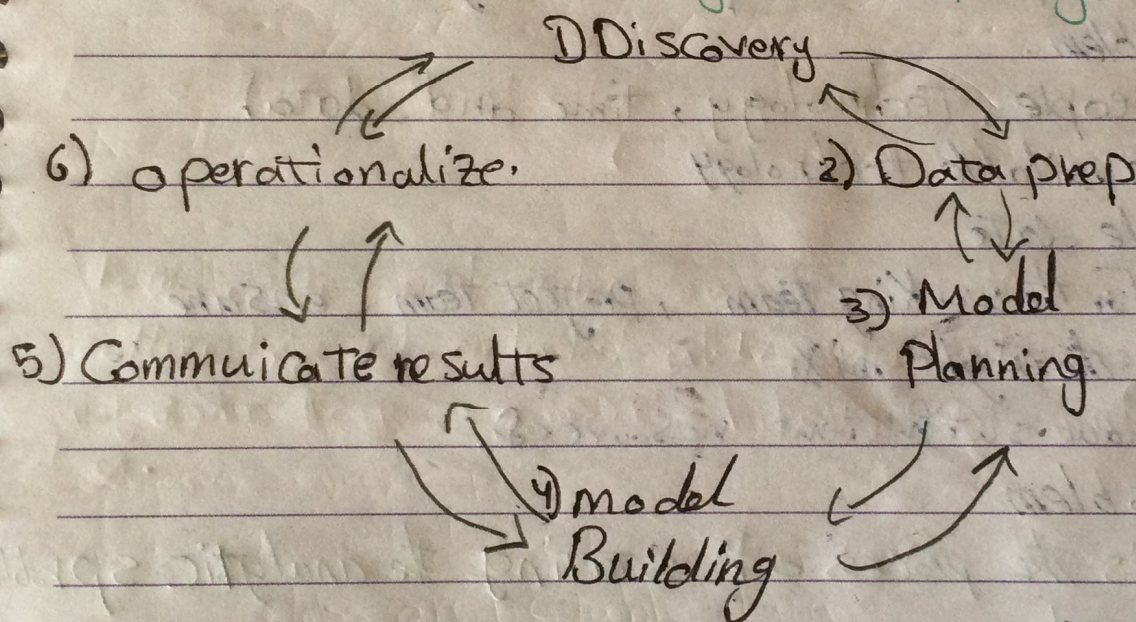
DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

**Database Administrator (DBA)** Database Administrator who provisions and configures database environment to support the analytical needs of working team

**Data Scientist** provide subject matter expertise for analytical techniques, data modelling, applying valid analytical techniques to given business problems and ensuring overall analytical objectives are met

## Data Analytics Life Cycle



## \*Data Analytics Life Cycle :-

We can go back and refine work done in prior phases given new insights and information that you've uncovered



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

## Phase 1: Discovery.

### • learn the Business Domain.

determine amount of domain knowledge needed to orient you to data and interpret results down stream.

### • Determine the general analytic problem type

« such as clustering, classification »

• then conduct initial research to learn about the domain area you'll be analyzing.

### • learn From the past.

• learn how previous attempts in the organization to solve this problem.

Resources (people, technology, time and data)

• assess available technology

• Available data

• people for the working team, project team ensure we have the right mix.

• Do you have sufficient resources?

## Frame the problem

Framing.  $\Rightarrow$  is the process of stating the analytic problem to be solved.

• state the analytics problem, why it is important

## Summary of discovery phase.

1) learn about domain knowledge.

2) detect available resources.

1) ... - technology - data - Tools - ...

« classification » clustering learning

((ALQSA)))



DATE: Dec 4

SUBJECT: \_\_\_\_\_

## "data preparation phase"

اعتبر أطول وأصعب مرحلة.

حيث هنا يكون هنا Sand Box وقد يكون في بعض الأحيان البريكتر من data warehouse.

1- بآلات data المتاحة هذه

2- تكون Sand Box هنا ينقل Limited data

ويوجد هنا 2 technologies

1) ETL Extract Transform Load.

2) ELT Extract local Transform.

حيث يتجمع ويتنقل البيانات المتاحة وهو يقرر إلى هو أين هو العمل "deep analysis"  
وهنا يتم التعامل مع DWH, it

Valid

Concurrence

process

missing Fields

check data. -x

1. Check data type (structured, unstructure, semi structured),

2. systematic error

1- prepare Analytic Sand Box

2- perform ELT

3- Familiarize yourself with the data thoroughly

4- Data Conditioning

5- Survey & visualize.

6- Determine methods

7- Techniques & work Flow.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

### \* phase 3: Model planning

- 1- Data exploration
- 2- variable selection
- 3- Model Selection

١- تحديد Framework بناء على الـ Technique التي ستستخدمها  
٢- Feature Selection, workflow التي ستستخدمها، نوع model  
٣- تحديد الشكل النهائي لـ model بناء على

### \* phase 4: model Building

- 1- develop data sets for testing, Training, and production purposes.
- 2- get the best environment you can for building models and work flows. R, RPL, ...

١- اختيار أنواع data وهل ستكون مدمجة ولا لا  
٢- التقييم في models هل أمكنه إنتاج شيء أو بعض الأخطاء  
٣- important parameters - valid - accuracy, mistakes

### \* phase 5: Communicate Results

- ١- أخذ محاولة اقناع sponsor بالقرارات التي تم كديرها سابقاً.  
٢- وعرض النتائج بـ clear.

Did we succeed? Did we Fail?



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

\* phase 6 so operationalize

• sub Test JAL system & check if it works

4 Core Deliverables to Meet most Stakeholder Needs

1 - presentation for project sponsors

- Big picture takeaways for executive level stakeholders.
- Determine key messages to aid their decision-making process
- Focus on clean, easy visuals for the presenter to explain and for the viewer to grasp.

2 - presentation for Analysis

- Business process changes
- reporting changes
- Follow Data Scientists will want the details and are comfortable with technical graphs

3 - Code for technical people.

4 - Technical specs of implementing the code.

⇒ Analyst wish list for a successful Analytics project.

\* Data & workspaces

- 1 - Access all data
- 2 - up-to-date data dictionary
- 3 - Ability to move data back between staging
- 4 - Sand box.



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

Tools:.

- statistical, mathematical, visual SW.
- Tool or place to log errors with systems.
- Collaboration → online platform for communication with team members.



DATE: lec

SUBJECT: \_\_\_\_\_

estimated	Actual	
+ve	+ve	→ TP
-ve	+ve	→ TN
+ve	-ve	→ FP
-ve	+ve	→ FN

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FN + FP}$$

- which is
- Building model that solve the problem
- Check accuracy of model → Compare with myself  
→ Compare with others Same problem

\* hypothesis testing. وذلك بنفس نتائج نفس Samples بناءً على ما هو

→ Choosing variable → لو عمل فرقاً في الاختيار  
لو نفس العمل فرقاً في النتائج عنه وذلك لعدم وجود Variance

الغرض من بناء model هنا ليس تأكيد نتائج ولكن من أجل تقييم نتائج وجوده بالفعل.

\* median. لو فردى ييقن المتصف  
فردى ينقسم 2 في النصف وتقسيم على 2.

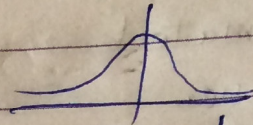
وإذا كان mean هناك احتمال كبير على 0  
مما قد يكون هناك 2 في بعض الأحيان قد تكون + و -  
في مختلفين لذلك تم اللجوء إلى variance



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

Test check difference of mean.



1. T-test  $\rightarrow$  normal distribution.  $\leftarrow$  الأكثر دقة
2. welch's T-test  $\rightarrow$  default to PR  $\leftarrow$  more tightly.
3. Rank sum  $\rightarrow$  general  $\leftarrow$  less tightly.

$H_0 \Rightarrow$  "hypothesis"  $\leftarrow$  وتعتبر من حيث أي شيء  
 $H_1 \Rightarrow$  " "  $\leftarrow$  تعني وجود فرق بين

T-test "two way test"  $\leftarrow$  inter variance  $\leftarrow$  توزيع  $\leftarrow$  distribution.

أما Student T-test  $\leftarrow$  variance  $\leftarrow$  variance

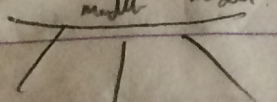
①  $t=0$  accept null hypothesis.  $\leftarrow$  من فرق  
 $t=w$  area under the curve.  $\leftarrow$  في فرق  
 accept inter hypothesis.  $\leftarrow$  بين فرق

Rank Sum.

$w \rightarrow$  old new threshold value.

Samples			A	
$x_1$	$x_2$	$x_3$	$x_{11}$	$x_{21}$
			$x_{12}$	$x_{22}$

$$w = \sum \text{sgn}(-\text{old} + \text{new})$$



$x \rightarrow 0$   $\leftarrow$  ((ALQSA))  
 $\text{sgn}(x) \rightarrow 0$   $\leftarrow$   $+$   $\leftarrow$   $-$   
 $\text{sgn}(x) \rightarrow 0$   $\leftarrow$   $1$   $\leftarrow$   $-1$



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

Power  $\Rightarrow$  Positive value of  $F$   
 Signif  $\Rightarrow$  F-P error rate  
 effect size  $\Rightarrow$  Actual difference between 2 means -

ANOVA

Two way  $\vee$  one way  $\rightarrow$  2 variables

new old  $\rightarrow$  2 model  $\rightarrow$  2 variables  
 2 model  $\rightarrow$  2 variables  $\rightarrow$  2 variables

5.28

1. Calculate mean of each population.

$$m_1 = 2.67$$

$$m_2 = 2.67$$

$$m_3 = 3$$

①	②	③
1	2	2
2	4	3
5	2	4

grand mean

$$\bar{m}_0 = \frac{m_1 + m_2 + m_3}{n} = \frac{2.67 + 2.67 + 3}{9} = 2.78$$

2. Sum of Squares (SS)

$$SS_{within} = \sum (X_i - m_i)^2 + \sum (X_2 - m_2)^2 + \sum (X_3 - m_3)^2$$

$$(1 - 2.67)^2 + (2 - 2.67)^2 + (5 - 2.67)^2 + (2 - 2.67)^2 + (4 - 2.67)^2$$

$$SS_{total} = \sum (X - m_0)^2 = 13.6$$

(((ALQSA)))



DATE: \_\_\_\_\_

SUBJECT: \_\_\_\_\_

•  $SS_{\text{Between}} = SS_{\text{total}} - SS_{\text{within}} = 0.23$

•  $\sigma_w^2 = V_w = \frac{SS_w}{N-K} = 13.34 / (9-3) = 2.22$

عدد النماذج  $K=3$  ←  $N-K$  عدد العينات  $N=9$

•  $\sigma_B^2 = \frac{SS_B}{K-1} = \frac{0.23}{2} = 0.12$

2) accept null hypothesis.

•  $F = \sigma_B^2 / \sigma_w^2 \rightarrow$

$\sigma_w$  → انحراف المعياري  
 Variance → تباين  
 أكبر

